

# Self-Stabilizing $k$ -out-of- $\ell$ Exclusion on Tree Networks

Ajoy K. Datta<sup>1</sup>   Stéphane Devismes<sup>2</sup>   Florian Horn<sup>3</sup>   Lawrence L. Larmore<sup>1</sup>

<sup>1</sup>School of Computer Science, University of Nevada Las Vegas, {datta,larmore}@cs.unlv.edu

<sup>2</sup>VERIMAG, Université Grenoble 1 Joseph Fourier, stephane.devismes@imag.fr

<sup>3</sup>LIAFA, Université Paris 7 Denis Diderot, horn@liafa.jussieu.fr

## Abstract

In this paper, we address the problem of  $k$ -out-of- $\ell$  exclusion, a generalization of the mutual exclusion problem, in which there are  $\ell$  units of a shared resource, and any process can request up to  $k$  units ( $1 \leq k \leq \ell$ ). We propose the first deterministic self-stabilizing distributed  $k$ -out-of- $\ell$  exclusion protocol in message-passing systems for asynchronous oriented tree networks which assumes bounded local memory for each process.

**Keywords:** Fault-tolerance, self-stabilization, resource allocation,  $k$ -out-of- $\ell$  exclusion, oriented tree networks.

## 1 Introduction

The basic problem in resource allocation is the management of shared resources, such as printers or shared variables. The use of such resources by an agent affects their availability for the other users. In the aforementioned cases, at most one agent can access the resource at any time, using a special section of code called a *critical section*. The associated protocols must guarantee the *mutual exclusion* property [13]: the critical section can be executed by at most one process at any time. The  $\ell$ -exclusion property [6] is a generalization of mutual exclusion, where  $\ell$  processes can execute the critical section simultaneously. Thus, in  $\ell$ -exclusion,  $\ell$  units of a same resource (*e.g.*, a pool of IP addresses) can be allocated. This problem can be generalized still further by considering heterogeneous requests, *e.g.*, bandwidth for audio or video streaming. The  $k$ -out-of- $\ell$  exclusion property [12] allows us to deal with such requests; requests may vary from 1 to  $k$  units of a given resource, where  $1 \leq k \leq \ell$ .

**Contributions.** In this paper, we propose a (deterministic) *self-stabilizing* distributed  $k$ -out-of- $\ell$  exclusion protocol for asynchronous oriented tree networks. A protocol is self-stabilizing [5] if, after transient faults hit the system and place it in some arbitrary global state, the systems recovers from this catastrophic situation without external (*e.g.* human) intervention in finite time. Our protocol is written in the message-passing model, and assumes bounded memory per process. To the best of our knowledge, there is no prior protocol of this type in the literature.

Obtaining a self-stabilizing solution for the  $k$ -out-of- $\ell$  exclusion problem in oriented trees is desirable, but also complex. Our main reason for dealing with oriented trees is that extension to general rooted networks is trivial; it consists of running the protocol concurrently with a spanning tree construction (for message passing systems), such as given in [1, 4]. In the other hand, the complexity of the solution comes from the fact that the problem is a generalization of mutual exclusion. This is exacerbated by the difficulty of obtaining self-stabilizing solutions in message-passing system (the more realistic model), as underlined by the impossibility result of Gouda and Multari [7].

Designing protocols for such problems on realistic systems often leads to obfuscated solutions. A direct consequence is then the difficulty of checking, or analyzing the solution. To circumvent this problem, we propose, here, a step-by-step approach. We start from a “naive” non-operating circulation of  $\ell$  resource tokens. Incrementally, we augment this solution with several other types of tokens until we obtain a correct non fault-tolerant solution. We then introduce an additional control mechanism that guarantees self-stabilization assuming unbounded local memory. Finally, we modify the protocol to accommodate bounded local memory.

We validate our approach by showing correctness and analyzing waiting time, a crucial parameter in resource allocation.

**Related Work.** Two kinds of protocols are widely used in the literature to solve the  $k$ -out-of- $\ell$  exclusion problem: permission-based protocols, and  $\ell$ -token circulation. All non self-stabilizing solutions currently in the literature are permission-based. In a permission-based protocol, any process can access a resource after receiving permissions from all processes [12], or from the processes constituting its quorum [10, 11]. There exist two self-stabilizing solutions for  $k$ -out-of- $\ell$  exclusion on the oriented rooted ring [2, 3]. These solutions are based on circulation of  $\ell$  tokens, where each token corresponds to a resource unit.

**Outline.** The remainder of the paper is organized as follows: In the next section, we define the model used in this paper. In Section 3, we present our self-stabilizing  $k$ -out-of- $\ell$  exclusion protocol. In Section 4, we provide the proof of correctness of our protocol, and we analyze its waiting time. Finally, we conclude in Section 5.

## 2 Preliminaries

**Distributed Systems.** We consider *asynchronous distributed systems* having a *finite* number of *processes*. By asynchronous, we mean that there is no bound on message delay, clock drift, or process execution rate. Every process can directly communicate with a subset of processes called *neighbors*. We denote by  $\Delta_p$  the number of neighbors of a process  $p$ . We consider the message-passing model where communication between neighboring processes is carried out by *messages* exchanged through *bidirectional links*, *i.e.*, each link can be seen as two channels in opposite directions. The neighbor relation defines a *network*. We assume that the topology of the network is that of an *oriented tree*. *Oriented* means that there is a distinguished process called *root* (denoted  $r$ ) and that every non-root process knows which neighbor is its *parent* in the tree, *i.e.*, the neighbor that is nearest to the root. We say that process  $q$  is a *child* of process  $p$  if and only if  $p$  is the parent of  $q$ .

A process is a sequential deterministic machine with input/output capabilities and bounded local memory, and that uses a local algorithm. Each process executes its local algorithm by taking *steps*. In a step, a process executes two actions in sequence: (1) either it tries to receive a message from another process, sends a message to another process, or does nothing; and then (2) modifies some of its variables.<sup>1</sup> The local algorithm is structured as infinite loop that contains finitely many actions.

We assume that the channels incident to a process  $p$  are locally distinguished by a *label*, a number in the range  $\{0 \dots \Delta_p - 1\}$ ; by an abuse of notation, we may refer to a neighbor  $q$  of  $p$  by the label of  $p$ 's channel to  $q$ . We assume that the channels are *reliable*, meaning that no message can be lost (after transient faults are corrected) and *FIFO*, meaning that messages are received in the order they are sent. We also assume that each channels initially contains some arbitrary messages, but not more than a given bound  $C_{\text{MAX}}$ .<sup>2</sup>

A message is of the following form:  $\langle \text{type}, \text{value} \rangle$ . The *value* field is omitted if the message does not carry any value. A message may also contain more than one value.

A *distributed protocol* is a collection of  $n$  local algorithms, one for each process. We define the *state* of each process to be the state of its local memory and the contents of its incoming channels. The global state of the system, referred to as a *configuration*, is defined as the product of the states of processes. We denote by  $\mathcal{C}$  the set of all possible configuration. An *execution* of a protocol  $\mathcal{P}$  in a system  $\mathcal{S}$  is an infinite sequence of configurations (of  $\mathcal{S}$ )  $\gamma_0 \gamma_1 \dots \gamma_i \dots$  such that in any transition  $\gamma_i \mapsto \gamma_{i+1}$  either a process take a step, or an external (*w.r.t.* the protocol) application modifies an input variable. Any execution is assumed to be *asynchronous* but *fair*: Every process takes an infinite number of steps in the execution but the time between two steps of a process is unbounded.

**$k$ -out-of- $\ell$  exclusion.** In the  $k$ -out-of- $\ell$  exclusion problem, the existence of  $\ell$  units of a shared resource is assumed. Any process can request at most  $k$  units of the shared resource, where  $k \leq \ell$ . We say that a protocol satisfies the  $k$ -out-of- $\ell$  exclusion specification if it satisfies the following three properties:

<sup>1</sup>When there is ambiguity, we denote by  $x_p$  the variable  $x$  in the code of process  $p$ .

<sup>2</sup>This assumption is required to obtain a deterministic self-stabilizing solution working with bounded process memory; see [7].

- **Safety:** At any given time, each resource unit (*n.b.*, here a resource unit corresponds to a token) is used by at most one process, each process uses at most  $k$  resource units, and at most  $\ell$  resource units are used.
- **Fairness:** If a process requests at most  $k$  resource units, then its request is eventually satisfied (*i.e.* it can eventually use the resource unit it requests using a special section of code called *critical section*).
- **Efficiency:** Informally, this means that as many requests as possible must be satisfied simultaneously.

The above mentioned notion of *efficiency* is difficult to define precisely. A convenient parameter was introduced in [3] to formally characterize efficiency:  $(k, \ell)$ -*liveness*, defined as follows. Assume that there is a subset  $I$  of processes such that every process in  $I$  is executing its critical section forever (*i.e.*, it holds some resource units forever). Let  $\alpha$  be the total number of resource units held forever by the processes in  $I$ . Let  $R$  be the set of processes not in  $I$  that are requesting some resource units; for each  $q \in R$ , let  $r_q$  be the number of resource units being requested by  $q$ , and assume that  $r_q \leq \ell - \alpha$  for all  $q \in R$ . Then, if  $R \neq \emptyset$ , at least one member of  $R$  eventually satisfies its request.

**Waiting Time.** The *waiting time* [14] is the maximum number of times that all processes can enter in the critical section before some process  $p$ , starting from the moment  $p$  requests the critical section.

**Interface.** In any  $k$ -out-of- $\ell$  exclusion protocol, a process needs to interact with the application that requests the resource units. To manage these interactions, we use the following interface at each process:

- **State**  $\in \{Req, In, Out\}$ . **State** = *Req* means that the application is requesting some resource units. **State** switches from *Req* to *In* when the application is allowed to access to the requested resource units. **State** switches from *In* to *Out* when the requested resource units are released into the system. The switching of **State** from *Req* to *In* and from *In* to *Out* is managed by the  $k$ -out-of- $\ell$  exclusion protocol itself; while the switching from *Out* to *In* is managed by the application. Other transitions (for instance, *In* to *Req*) are forbidden.
- **Need**  $\in \{0 \dots k\}$ , the number of resource units currently being requested by the application.
- **EnterCS()**: *function*. This function is called by the protocol to allow the application to execute the *critical section*. From this call, the application has control of the resource units until the end of the critical section (we assume that the critical section is always executed in finite, yet unbounded, time).
- **ReleaseCS()**: *Boolean*. This predicate holds if and only if the application is not executing its critical section.

**Self-Stabilization** [5]. A *specification* is a predicate over the set of all executions. A set of configurations  $\mathcal{C}_1 \subseteq \mathcal{C}$  is an *attractor* for a set of configurations  $\mathcal{C}_2 \subseteq \mathcal{C}$  if for any  $\gamma \in \mathcal{C}_2$  and any execution whose initial configuration is  $\gamma$ , the execution contains a configuration of  $\mathcal{C}_1$ .

**Definition 1** A protocol  $\mathcal{P}$  is self-stabilizing for the specification  $SP$  in a system  $\mathcal{S}$  if there exists a non-empty subset of  $\mathcal{L}$  such that:

- Any execution of  $\mathcal{P}$  in  $\mathcal{S}$  starting from a configuration of  $\mathcal{L}$  satisfies  $SP$  (Closure Property).
- $\mathcal{L}$  is an attractor for  $\mathcal{C}$  (Convergence Property).

### 3 Protocol

In this section we present our self-stabilizing  $k$ -out-of- $\ell$  exclusion protocol for oriented trees (Algorithms 1 and 2). Our solution uses circulation of several types of tokens. To clearly understand the function of these tokens, we adopt a step-by-step approach: we start from “naive” non-operating circulation of  $\ell$  resource tokens. Incrementally, we augment this solution with several other types of tokens, until we obtain a non-fault-tolerant solution. We then add an additional control mechanism that guarantees self-stabilization, assuming unbounded local memory of processes. Finally, we modify our protocol to work with bounded memory.

---

**Algorithm 1**  $k$ -out-of- $\ell$  exclusion on oriented trees, code for the root  $r$ 


---

```

1: variables:
2:    $C, \text{myC} \in [0 \dots 2(n-1)(C_{\max} + 1)]$ ;  $\text{Succ} \in [0 \dots \Delta_r - 1]$ 
3:    $\text{RSet}$ : multiset of at most  $k$  values taken in  $[0 \dots \Delta_r - 1]$ 
4:    $\text{Need} \in [0 \dots k]$ ;  $\text{State} \in \{\text{Req}, \text{In}, \text{Out}\}$ 
5:    $\text{Prio} \in \{\perp, 0, \dots, \Delta_r - 1\}$ 
6:    $R, \text{Reset}$ : Booleans;  $\text{SToken}, PT \in [0 \dots \ell + 1]$ 
7:    $\text{SPush}, \text{SPrio}, PPr \in [0 \dots 2]$ 
8: repeat forever
9:   for all  $q \in [0 \dots \Delta_r - 1]$  do
10:    if  $(\text{receive}(\text{ResT}) \text{ from } q) \wedge \neg \text{Reset}$  then
11:      if  $(\text{State} = \text{Req}) \wedge (|\text{RSet}| < \text{Need})$  then
12:         $\text{RSet} \leftarrow \text{RSet} \cup \{q\}$ 
13:      else
14:        if  $q = \Delta_r - 1$  then
15:           $\text{SToken} \leftarrow \min(\text{SToken} + 1, \ell + 1)$ 
16:        end if
17:         $\text{send}(\text{ResT}) \text{ to } q + 1$ 
18:      end if
19:    end if
20:    if  $(\text{receive}(\text{PushT}) \text{ from } q) \wedge \neg \text{Reset}$  then
21:      if  $(\text{Prio} \neq \perp) \wedge (\text{State} \neq \text{Req} \vee |\text{RSet}| < \text{Need}) \wedge$ 
22:         $(\text{State} \neq \text{In})$  then
23:        for all  $i \in \text{RSet}$  do
24:          if  $i = \Delta_r - 1$  then
25:             $\text{SToken} \leftarrow \min(\text{SToken} + 1, \ell + 1)$ 
26:          end if
27:           $\text{send}(\text{ResT}) \text{ to } i + 1$ 
28:        end for
29:         $\text{RSet} \leftarrow \emptyset$ 
30:      end if
31:      if  $q = \Delta_r - 1$  then
32:         $\text{SPush} \leftarrow \min(\text{SPush} + 1, 2)$ 
33:      end if
34:       $\text{send}(\text{PushT}) \text{ to } q + 1$ 
35:    end if
36:    if  $(\text{receive}(\text{PrioT}) \text{ from } q) \wedge \neg \text{Reset}$  then
37:      if  $\text{Prio} = \perp$  then
38:         $\text{Prio} \leftarrow q$ 
39:      else
40:         $\text{send}(\text{PrioT}) \text{ to } q + 1$ 
41:      end if
42:    end if
43:    if  $(\text{receive}(\text{ctrl}, C, R, PT, PPr) \text{ from } q)$  then
44:      if  $(q = \text{Succ}) \wedge (\text{myC} = C)$  then
45:         $\text{Succ} \leftarrow \text{Succ} + 1$ 
46:        if  $\text{Succ} = 0$  then
47:           $\text{myC} \leftarrow \text{myC} + 1$ 
48:           $\text{Reset} \leftarrow (PT + \text{SToken} > \ell) \vee$ 
49:             $(PPr + \text{SPrio} > 1) \vee (\text{SPush} > 1)$ 
50:          if  $\text{Reset}$  then
51:             $\text{RSet} \leftarrow \emptyset$ 
52:             $\text{Prio} \leftarrow \perp$ 
53:          end if
54:        else
55:          if  $PPr + \text{SPrio} < 1$  then
56:             $\text{send}(\text{PrioT}) \text{ to } 0$ 
57:          end if
58:          while  $PT + \text{SToken} < \ell$  do
59:             $\text{send}(\text{ResT}) \text{ to } 0$ 
60:             $\text{SToken} \leftarrow \min(\text{SToken} + 1, \ell + 1)$ 
61:          end while
62:          if  $\text{SPush} < 1$  then
63:             $\text{send}(\text{PushT}) \text{ to } 0$ 
64:          end if
65:          end if
66:           $\text{SToken} \leftarrow 0$ 
67:           $\text{SPrio} \leftarrow 0$ 
68:           $\text{SPush} \leftarrow 0$ 
69:           $PT \leftarrow 0$ 
70:           $PPr \leftarrow 0$ 
71:        end if
72:         $PT \leftarrow \min(PT + |\text{RSet}|_q, \ell + 1)$ 
73:        if  $\text{Prio} = q$  then
74:           $PPr \leftarrow \min(PPr + 1, 2)$ 
75:        end if
76:         $\text{send}(\text{ctrl}, \text{myC}, \text{Reset}, PT, PPr) \text{ to Succ}$ 
77:         $\text{RestartTimer}()$ 
78:      end if
79:    end if
80:    if  $(\text{State} = \text{Req}) \wedge (|\text{RSet}| \geq \text{Need})$  then
81:       $\text{State} \leftarrow \text{In}$ 
82:       $\text{EnterCS}()$ 
83:    end if
84:    if  $(\text{State} = \text{In}) \wedge \text{ReleaseCS}()$  then
85:      for all  $i \in \text{RSet}$  do
86:        if  $i = \Delta_r - 1$  then
87:           $\text{SToken} \leftarrow \min(\text{SToken} + 1, \ell + 1)$ 
88:        end if
89:         $\text{send}(\text{ResT}) \text{ to } i + 1$ 
90:      end for
91:       $\text{RSet} \leftarrow \emptyset$ 
92:       $\text{State} \leftarrow \text{Out}$ 
93:    end if
94:    if  $(\text{Prio} \neq \perp) \wedge (\text{State} \neq \text{Req} \vee |\text{RSet}| \geq \text{Need})$  then
95:      if  $\text{Prio} = \Delta_r - 1$  then
96:         $\text{SPrio} \leftarrow \min(\text{SPrio} + 1, 2)$ 
97:      end if
98:       $\text{send}(\text{PrioT}) \text{ to Prio} + 1$ 
99:       $\text{Prio} \leftarrow \perp$ 
100:    end if
101:    if  $\text{TimeOut}()$  then
102:       $\text{send}(\text{ctrl}, \text{myC}, \text{Reset}, 0, 0) \text{ to Succ}$ 
103:       $\text{RestartTimer}()$ 
104:    end if
105:  end repeat

```

---

---

**Algorithm 2**  $k$ -out-of- $\ell$  exclusion on oriented trees, code for the other process  $p$ 


---

```

1: variables:
2:  $C, \text{myC} \in [0 \dots 2(n-1)(C_{\text{MAX}} + 1)]$ ;  $\text{Succ} \in [0 \dots \Delta_p - 1]$ 
3:  $\text{RSet}$ : multiset of at most  $k$  values taken in  $[0 \dots \Delta_p - 1]$ 
4:  $\text{Need} \in [0 \dots k]$ ;  $\text{State} \in \{\text{Req}, \text{In}, \text{Out}\}$ 
5:  $\text{Prio} \in \{\perp, 0, \dots, \Delta_p - 1\}$ 
6:  $R, \text{Ok}$ : Booleans;  $PT \in [0 \dots \ell + 1]$ ;  $PPr \in [0 \dots 2]$ 
7: repeat forever
8:   for all  $q \in [0 \dots \Delta_p - 1]$  do
9:     if (receive( $\text{ResT}$ ) from  $q$ ) then
10:       if ( $\text{State} = \text{Req}$ )  $\wedge$  ( $|\text{RSet}| < \text{Need}$ ) then
11:          $\text{RSet} \leftarrow \text{RSet} \cup \{q\}$ 
12:       else
13:         send( $\text{ResT}$ ) to  $q + 1$ 
14:       end if
15:     end if
16:     if (receive( $\text{PushT}$ ) from  $q$ ) then
17:       if ( $\text{Prio} \neq \perp$ )  $\wedge$  ( $\text{State} \neq \text{Req} \vee |\text{RSet}| < \text{Need}$ )  $\wedge$ 
18:         ( $\text{State} \neq \text{In}$ ) then
19:         for all  $i \in \text{RSet}$  do
20:           send( $\text{ResT}$ ) to  $i + 1$ 
21:         end for
22:          $\text{RSet} \leftarrow \emptyset$ 
23:       end if
24:       send( $\text{PushT}$ ) to  $q + 1$ 
25:     end if
26:     if (receive( $\text{PrioT}$ ) from  $q$ ) then
27:       if  $\text{Prio} = \perp$  then
28:          $\text{Prio} \leftarrow q$ 
29:       else
30:         send( $\text{PrioT}$ ) to  $q + 1$ 
31:       end if
32:     end if
33:     if (receive( $\text{ctrl}, C, R, PT, PPr$ ) from  $q$ ) then
34:        $\text{Ok} \leftarrow \text{false}$ 
35:       if ( $q = \text{Succ}$ )  $\wedge$  ( $\text{myC} = C$ )  $\wedge$  ( $\text{Succ} \neq 0$ ) then
36:          $\text{Succ} \leftarrow \text{Succ} + 1$ 
37:          $\text{Ok} \leftarrow \text{true}$ 
38:         if  $R$  then
39:            $\text{RSet} \leftarrow \perp$ 
40:         end if
41:       end if
42:       if ( $q = 0$ ) then
43:          $\text{Ok} \leftarrow \text{true}$ 
44:         if  $\text{myC} \neq C$  then
45:            $\text{Succ} \leftarrow \min(1, \Delta_p - 1)$ 
46:           if  $R$  then
47:              $\text{RSet} \leftarrow \emptyset$ 
48:              $\text{Prio} \leftarrow \perp$ 
49:           end if
50:         end if
51:          $\text{myC} \leftarrow C$ 
52:       end if
53:       if  $\text{Ok}$  then
54:          $PT \leftarrow \min(PT + |\text{RSet}|_q, \ell + 1)$ 
55:         if  $\text{Prio} = q$  then
56:            $PPr \leftarrow \min(PPr + 1, 2)$ 
57:         end if
58:         send( $\text{ctrl}, \text{myC}, R, PT, PPr$ ) to  $\text{Succ}$ 
59:       end if
60:     end if
61:   end for
62:   if ( $\text{State} = \text{Req}$ )  $\wedge$  ( $|\text{RSet}| \geq \text{Need}$ ) then
63:      $\text{State} \leftarrow \text{In}$ 
64:      $\text{EnterCS}()$ 
65:   end if
66:   if ( $\text{State} = \text{In}$ )  $\wedge$   $\text{ReleaseCS}()$  then
67:     for all  $i \in \text{RSet}$  do
68:       send( $\text{ResT}$ ) to  $i + 1$ 
69:     end for
70:      $\text{RSet} \leftarrow \emptyset$ 
71:      $\text{State} \leftarrow \text{Out}$ 
72:   end if
73:   if ( $\text{Prio} \neq \perp$ )  $\wedge$  ( $\text{State} \neq \text{Req} \vee |\text{RSet}| \geq \text{Need}$ ) then
74:     send( $\text{PrioT}$ ) to  $\text{Prio} + 1$ 
75:      $\text{Prio} \leftarrow \perp$ 
76:   end if
77: end repeat

```

---

**A non-fault-tolerant protocol.** The basic principle of our protocol is to use  $\ell$  circulating *resource tokens* (the  $\text{ResT}$  messages) following depth-first search (DFS) order: when a process  $p$  receives a token from channel number  $i$ , and if that token is retransmitted, either immediately or later, it will be sent to its neighbor along channel number  $i + 1$  (modulo  $\Delta_p$ ). (This same rule will also be followed by all the types of tokens we will later describe.) Figure 1 shows the path followed by a token during depth-first circulation in an oriented tree (recall that any non-root process locally numbers the channel to its parent by 0). In this way, the oriented tree emulates a ring with a designated leader (see Figure 4), and we refer to the path followed by the tokens as the *virtual ring*.

As explained Section 2, the requests are managed by the variables **State** and **Need**. Each process also uses the multiset<sup>3</sup> variable **RSet** to collect the tokens; the collected tokens are said to be “reserved.” While  $\text{State} = \text{Req}$  and  $|\text{RSet}| < \text{Need}$ , a process collects all tokens it receives; it also stores in **RSet** the number of the channel from which it receives each token, so that when it is finally retransmitted, it will continue its correct path around the virtual ring. When  $\text{State} = \text{Req}$  and  $|\text{RSet}| \geq \text{Need}$ , it enters the critical section: **State** is set to *In* and the function  $\text{EnterCS}()$  is called. Once the critical section is done (*i.e.*, when  $\text{State} = \text{In}$  and the predicate  $\text{ReleaseCS}()$  holds) **State** is set to *Out*, all tokens in **RSet** are retransmitted, and **RSet** is set to  $\emptyset$ . When a process receives a token it does not need, it immediately retransmits it.

Unfortunately, such a simple protocol does not always guarantee liveness. Figure 2 shows a case where liveness is not maintained. In this example, there are five resources tokens (*i.e.*,  $\ell = 5$ ) and each process can request up to three tokens (*i.e.*,  $k = 3$ ). In the configuration shown on the left side of the figure, processes  $a$ ,  $b$ ,  $c$ , and  $d$  request more tokens than they will receive. This configuration will lead to the deadlock configuration shown on the right side of the figure: processes  $a$ ,  $b$ ,  $c$ , and  $d$  reserve all the tokens they receive and never release them because their requests are never satisfied.

---

<sup>3</sup>*N.b.* a multiset can contain several identical items.

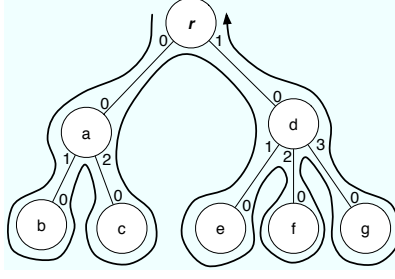


Figure 1: Depth-first token circulation on oriented trees.

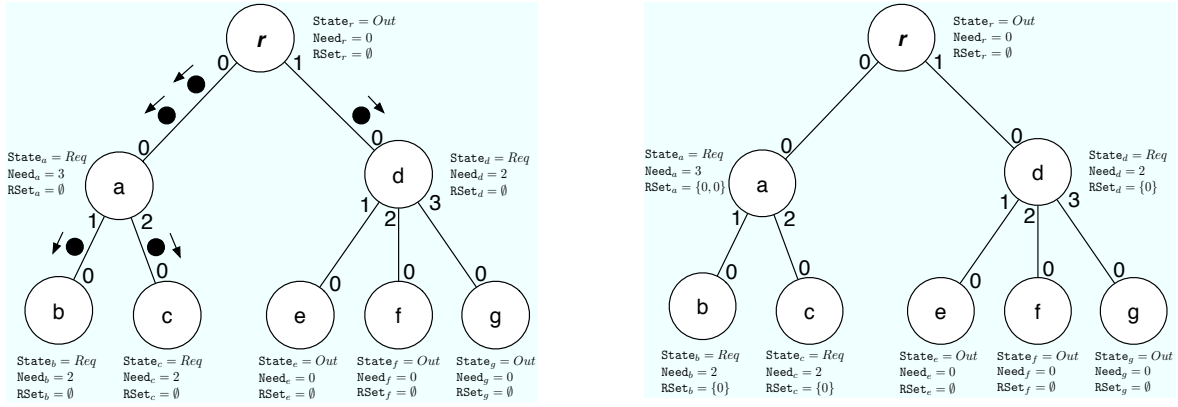


Figure 2: Possible deadlock.

We can prevent deadlock by adding a new type of token, called the *pusher* (the message `PushT`). If the system is in a legitimate state, there is exactly one pusher. It permanently circulates through the virtual ring, and prevents a process that is not in the critical section from holding resource tokens forever. When a process receives the pusher, it releases all its reserved tokens, unless if it is either in its critical section ( $State = In$ ) or is enabled to enter its critical section ( $State = Req$  and  $|RSet| \geq Need$ ). In either case, it retransmits the pusher.

The pusher protects the system from deadlock. However, it can cause *livelock*; an example is shown in Figure 3, for 2-out-of-3 exclusion in a tree of three processes. In Configuration (i), every process is a requester:  $r$  and  $b$  request one resource token and  $a$  requests two resource tokens. Also, every process has a resource token in one of its incoming channels, and none holds any resource token. Finally, the pusher is in the channel from  $a$  to  $r$  behind a resource token. Every process will collect the incoming resource token, and the system will reach the Configuration (ii) where  $r$  and  $b$  execute their critical section while  $a$  is still waiting for a resource token and the pusher is reaching  $r$ . When  $r$  receives the pusher, it retransmits it to  $b$ , while keeping its resource token, as shown in Configuration (iii). Similarly,  $b$  receives the pusher while executing its critical section, and retransmits it immediately to  $r$ , as shown in Configuration (iv), after which  $r$  retransmits the pusher to  $a$  (Configuration (v)). Assume now that  $a$  receives the pusher while  $r$  and  $b$  leave their critical sections. We obtain Configuration (vi):  $a$  must release its resource tokens because of the pusher. In Configuration (vii),  $r$  directly retransmits the resource token it receives because it is not a requester. Finally,  $r$  and  $b$  again become requesters for one resource token in Configuration (viii), which is identical to Configuration (i). We can repeat this cycle indefinitely, and process  $a$  never satisfies its request.

To solve this problem, we add a *priority token* (message `PrioT`) whose goal is to cancel the effect of the pusher. If the system is in a legitimate state, there is exactly one priority token. A process which receives the priority token retransmits it immediately, unless it has an unsatisfied request. In this case, the process holds the priority token (the variable `Prio` is set from  $\perp$  to the channel number from which the process receives

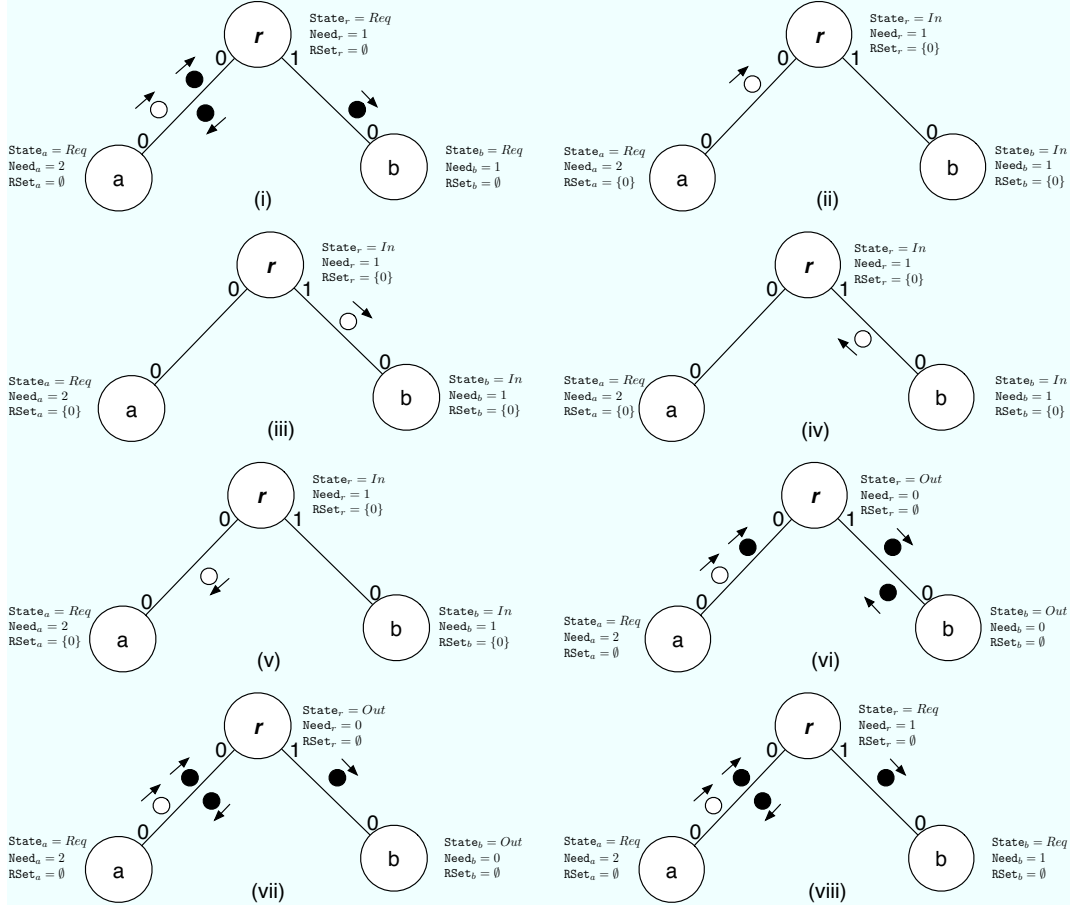


Figure 3: Possible livelock.

the priority token) until its request is satisfied: the token will then be released when the process enters its critical section. A process that holds the priority token does not release its reserved resource tokens when it receives the pusher: it only retransmits the pusher. As we will show later, this guarantees that the process will eventually satisfy its request.

Using these three types of tokens, we obtain a simple non self-stabilizing  $k$ -out-of- $\ell$  exclusion protocol. To make it self-stabilizing, we need additional structure.

**A controller for self-stabilization.** To achieve self-stabilization, we introduce one more type of token, the *controller*.

After a finite period of transient faults, some tokens may have disappeared or may be duplicated. To restore correct behavior, we need an additional self-stabilizing mechanism that regulates the number of tokens in the network: to achieve that, we use a mechanism similar to that introduced in [8] for self-stabilizing  $\ell$ -exclusion protocol on a ring. This mechanism is based on snapshot/reset technique.

The controller is a special token (message `ctrl`) that counts the other tokens; when it returns to the root after one full circulation, the root learns the number of tokens of each type (resource, pusher, priority), and then adjusts these numbers as necessary.

The controller can also be effected by transient faults. We use Varghese’s *counter flushing* [15] technique to enforce *depth first token circulation* (DFTC) in the tree.

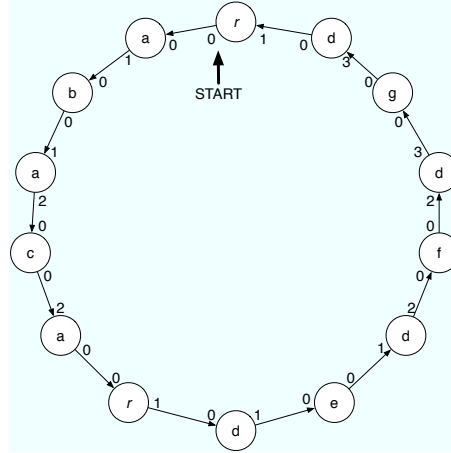


Figure 4: Virtual Ring.

We now explain how the resource tokens are counted by the controller. (It counts the other types of tokens similarly.) We split the count of the resource tokens into two subcounts:

- *The “passed” tokens.* When a process holds some resource tokens that came from channel  $i$  and receives the controller from the channel  $i$ , it retransmits the controller through channel  $i + 1$  while keeping the resource tokens: in this case, we say that the controller *passes* these tokens in the virtual ring. Indeed, these tokens were ahead the controller (in the virtual ring) before the process received the controller, and are behind afterward. The field *PT* of the controller message is used to compute the number of the passed resource tokens.
- *The tokens that are never passed by the controller.* These tokens are counted in the variable **SToken** maintained at the root. At the beginning of any circulation of the controller, the variable **SToken** is reset to 0. Then, until the end of the circulation of the controller, each time a resource token starts a new circulation (*i.e.* the token leaves the root from channel 0), **SToken** is incremented.

When the controller terminates its circulation, the number of resource tokens in the network is equal to  $PT + \text{SToken}$ , and the numbers of pusher tokens and priority tokens is likewise known to the root. Three cases are then possible:



- *The number of tokens is correct*, that is, there are  $\ell$  resource tokens, one pusher token, and one priority token. In this case, the system is stabilized.
- *There are too few tokens*. In this case, the root creates the number of additional tokens needed at the end of the traversal; the system is then stabilized.
- *There are too many tokens of some type*. In this case, we reset the network. We mark the controller token with a special flag (the field  $R$  in the message `ctrl`). The root transmits the marked controller, erases its reserved tokens as well as all the tokens it receives until the termination of the controller's traversal. Upon receiving the controller, every other process erases its reserved tokens. When the controller finishes its traversal, there is no token in the network. The root creates exactly  $\ell$  resource tokens, one pusher token and one priority token; and we are done.

**Self-stabilizing DFTC.** Using the counter flushing technique, we design a self-stabilizing *DFTC* to implement the controller. The principle of counter flushing is the following: after transient faults, the token message can be lost. Hence, the root must use a timeout mechanism to retransmit the token in case of deadlock. The timeout is managed using the function `RestartTimer()` (that allows it to reinitialize the timeout) and the predicate `TimeOut()` (which holds when a specified time interval is exceeded).<sup>4</sup>

Due to the use of the timeout, we must now deal with duplicated messages. Furthermore, arbitrary messages may exist in the network after faults (however they are assumed to be bounded). To distinguish the duplicates from the valid controller and to flush the system of corrupted messages, every process maintains a counter variable `myC` that takes values in  $\{0 \dots 2(n-1)(C_{\text{MAX}} + 1)\}$ , and marks each message with that value. Every process also maintains a pointer `Succ` to indicate to which process it must send the token. The effects of the reception of a token message differs for the root and the other processes:

- The root considers a token message as valid when the message comes from `Succ` and is marked with a value  $c$  such that `myC` =  $c$ . Otherwise, it simply ignores the message, meaning it does not retransmit it. If it receives a valid message, the root increments `Succ` (modulo  $\Delta_r$ ) and retransmits the token with the flag value `myC` to `Succ` so that the valid token follows DFS order. If `Succ` = 0, this means that the token just finished its previous circulation. As a consequence, the root increments `myC` (modulo  $2(n-1)(C_{\text{MAX}} + 1)$ ) before retransmitting the token.
- A non-root process  $p$  considers a message as valid in two cases: (1) When it receives a token message from its parent (channel 0) marked with a value  $c$  such that `myC`  $\neq c$  or (2) when it receives a token message from `Succ` and the message is marked with a value  $c$  such that `myC` =  $c$ . In case (1),  $p$  sets `myC` to  $c$  and `Succ` to  $\min(1, \Delta_p - 1)$  (*n.b.* in case of a leaf process `Succ` is set to 0) before retransmitting the token message marked with `myC` to `Succ`. In case (2),  $p$  increments `Succ` (modulo  $\Delta_p$ ) and then sends the token marked with `myC` to `Succ` so that the valid token follows DFS order. In all other cases,  $p$  considers the message to be invalid. In the case of an invalid message coming from channel 0 with `myC` =  $c$ ,  $p$  does not consider the message in the computation, but retransmits it to prevent deadlock. In all other cases,  $p$  simply ignores the message.

Using this method, the root increments its counter `myC` infinitely often and, due to the size of the `myC`'s domain, the `myC` variable of the root eventually takes a value that does not exist anywhere else in the system (because the number of possible values initially in the system is bounded by  $2(n-1)(C_{\text{MAX}} + 1)$ ). In this case, the token marked with the new value will be considered to be a valid token by every process. Until the end of that traversal, the root will ignore all other token messages. At the end of the traversal, the system will be stabilized.

**Dealing with bounded memory.** Due to the use of reset, the root does not need to know the exact number of tokens at the end of the controller's traversals. Actually, the root must only know if the number of tokens is too high, or the number of tokens it needs to add if the number is too low. Hence, the counting variables can be bounded by  $\ell + 1$  for the resource tokens and by 2 for the other types of token. The fact that a variable is assigned to its maximum value will mean that there are too many tokens in the network and so a reset must be started. Otherwise, the value of the counting variable will state whether there is a

---

<sup>4</sup>We assume that this time interval is sufficiently large to prevent congestion.

deficient number of tokens, and in that case, how many must be added. For any assignment to one of these bounded variables, the value is set to the minimum between its new computed value and the maximum value of its domain.

## 4 Correctness and Waiting Time

In this section, we first prove that our protocol is a self-stabilizing  $k$ -out-of- $\ell$  exclusion protocol. We then analyze its waiting time.

**Correctness.** We split the proof into three steps. (1) Recall that the controller part of our protocol is a self-stabilizing DFS token circulation. (2) We show that once the controller is stabilized to DFS token circulation, the system eventually stabilizes to the expected number of different tokens. (3) We show that once the system contains the expected number of tokens, the system stabilizes to the  $k$ -out-of- $\ell$  exclusion specification.

In our protocol, when a process receives a `ctrl` message, either it considers the message as valid or not. The process takes account of the messages for computations only when they are valid. Assume that a process  $p$  receives a `ctrl` message marked with the flag value  $c$  from channel  $q$ . Process  $p$  considers this message as valid if and only if  $(q = \text{Succ}_p \wedge c = \text{myC}_p) \vee (p \neq r \wedge (q = 0 \wedge c \neq \text{myC}_p))$ .

In the following, we call any `ctrl` message a *control token*. Each time a process receives a valid `ctrl` message, it makes some local computations, and then sends another `ctrl` message. In the case of a non-root process, the sent message is marked with the same flag as the received message: we consider it to be the same control token. In the case of the root, the sent message is marked either with the same value or with a new one. In the former case, we consider it to be the same control token, while in the latter case, we consider the received control token to have terminated its traversal, and the transmitted control token to be new.

To implement the control part, we use the counter flushing techniques introduced by Varghese in [15]. Hence, from [15], we can deduce the following lemma:

**Lemma 1** *Starting from any configuration, the system converges to a configuration at which:*

1. *There exists at most one valid control token in the network.*
2. *The root regularly creates a new valid control token.*
3. *Any valid control token visits all processes in DFS order.*

**Remark 1** *In our protocol, only the valid control tokens are considered in the computations. Hence, from now on, we only consider the valid control tokens and we simply refer to them as control tokens.*

We now show that starting from any configuration, the system eventually contains the expected number of each type of token.

Note that each resource token is either in a link (in this case, the token is said to be *free*) or it is stored in the `RSet` of a process (in this case, the token is said to be *reserved*). Hence, at any time the number of resource tokens in the network is equal to the sum of the sizes of the `RSet` multisets plus the number of *free* resource tokens.

Similarly, at any time, the number of priority tokens is equal to the number of processes satisfying `Prio`  $\neq \perp$  plus the number of *free* priority tokens.

Finally, as a process cannot store any pusher token, the number of pusher tokens is equal to the number of free pusher tokens.

**Lemma 2** *Let  $\gamma$  be the first configuration after the control part is stabilized. If, after  $\gamma$ , the root creates a control token whose reset field  $R$  is true, then the system contains no resource, priority, and pusher token at the end of the traversal of the control token.*

**Proof.** Consider any control token created by the root after configuration  $\gamma$ . Assume that the reset field  $R$  of the control token is set to true. Then, the `Reset` variable of the root is also true (see Line 73 in Algorithm 1). `Resetr` remains true until the control token terminates its traversal. Hence, during the traversal, any

token (except the control token) that is received by the root is ignored by the root and so disappears from the network (see Lines 10, 20, and 35 in Algorithm 1). Also, during its traversal, each process erases all tokens (except the control token) it holds when visited by the control token (see Line 48 in Algorithm 1, and Lines 37, and 46 in Algorithm 2). Hence, every resource, priority, or pusher token is either erased at a process when the process is visited by the control token, or is pushed to the root and then disappears. At the end of the traversal of the control token, the system contains no resource, priority, or pusher tokens.  $\square$

**Lemma 3** *Let  $\gamma$  be the first configuration after the control part is stabilized. When a control token created by the root after  $\gamma$  terminates its traversal, we have:*

- If  $PT + \text{SToken}_r > \ell$ , then there are more than  $\ell$  resource tokens in the network.
- If  $PT + \text{SToken}_r \leq \ell$ , then there are exactly  $PT + \text{SToken}_r$  resource tokens in the network.

**Proof.** Consider any control token created by the root after configuration  $\gamma$ . There are two cases:

- *The reset field  $R$  of the control token is true.* By Lemma 2, there is no resource token in the network when the control token terminates its circulation. So, to prove the lemma in this case, we must show that  $PT + \text{SToken}_r = 0$  at the end of the circulation.

First,  $\text{SToken}_r$  is reset to 0 (Line 63) before the control token starts its circulation (Line 73). Also,  $\text{Reset}_r$  is true when the control token starts its circulation (see Line 73 in Algorithm 1). Thus, until termination of the circulation,  $r$  ignores any resource tokens it receives (see Lines 10, 20, and 35) and so  $\text{SToken}_r$  is still equal to 0 at the end of the control token circulation.

Consider now the  $PT$  field of the control token. Before the start of the control token circulation,  $r$  executes the following action:  $\text{RSet}$  is set to  $\emptyset$  (Line 49 in Algorithm 1),  $PT$  is set to 0 (Line 66 in Algorithm 1), and, as a consequence,  $PT$  is set to  $\min(0, \ell + 1)$  (Line 69 in Algorithm 1). So, at the start of the control token circulation, the control token is sent with its  $PT$  field equal to 0. Since the reset field  $R$  of the control token is true, each time the control token arrives at a process, the process resets its  $\text{RSet}$  variable to  $\emptyset$  (see Lines 49 in Algorithm 1, Lines 38, and 47 in Algorithm 2) before setting  $PT$  to  $\min(PT + |\text{RSet}|_q, \ell + 1)$  (see Line 69 in Algorithm 1 and Line 54 in Algorithm 2) and then retransmitting the token. Hence,  $PT$  remains equal to 0 until the end of the circulation.

When the control token terminates its circulation,  $PT + \text{SToken}_r = 0$ , and we are done.

- *The reset field  $R$  of the control token is false.* In this case, we can remark that no resource token is erased during the circulation of the control token, because both  $R$  and  $\text{Reset}_r$  are false.

(\*) We now show that any resource token is counted at most once during the circulation of the control token. Due to the FIFO quality of the links and the fact that when the control token is received by a process, the process receives no other message before retransmitting the control token, we have the following property: a resource token is passed by the control token at most once during a circulation. So, during the circulation, either the resource token is counted into the  $PT$  field of the control token when the resource token is passed by the control token (see Line 69 in Algorithm 1 and Line 54 in Algorithm 2) or it is counted at the root when it terminates a loop of the virtual ring (Line 15). Hence, any resource token is counted at most once.

(\*\*) Finally we show, by contradiction, that any resource token is counted at least once during the circulation of the control token. Assume that a resource token is not counted during that circulation. Then, the resource token is never passed by the control token. The links are FIFO, and when the control token is received by a process, the process receives no other message before retransmitting the control token. So, the resource token is always ahead the control token in the virtual ring. As a consequence, the resource token is eventually counted at the root when it terminates a loop of the virtual ring (Line 15), contradiction.

From (\*), we know that if  $PT + \text{SToken}_r > \ell$  at the end of the control token circulation, then there are more than  $\ell$  resource tokens in the network. From (\*) and (\*\*), we know that if  $PT + \text{SToken}_r \leq \ell$  at the end of the control token circulation, then there are exactly  $PT + \text{SToken}_r$  resource tokens in the network, and we are done.

$\square$

Following similar reasoning, we obtain the following two lemmas:

**Lemma 4** *Let  $\gamma$  be the first configuration after the control part is stabilized. When a control token created by the root after  $\gamma$  terminates its traversal, we have:*

- *If  $\text{SPrio}_r + PPr > 1$ , there is more than one priority token in the network.*
- *If  $\text{SPrio}_r + PPr \leq 1$ , there are exactly  $\text{SPrio}_r + PPr$  priority tokens in the network.*

**Lemma 5** *Let  $\gamma$  be the first configuration after the control part is stabilized. When a control token created by the root after  $\gamma$  terminates its traversal, we have:*

- *If  $\text{SPush}_r > 1$ , then there is more than one pusher token in the network.*
- *If  $\text{SPush}_r \leq 1$ , then there are exactly  $\text{SPush}_r$  pusher tokens in the network.*

**Lemma 6** *Starting from any configuration, the system eventually reaches a configuration from which there always exist exactly  $\ell$  resource tokens.*

**Proof.** Let  $\gamma$  be the first configuration after the control part is stabilized. Consider any control token created by the root after  $\gamma$ . Let us study the two following cases:

- *$PT + \text{SToken}_r \leq \ell$  at the end of the control token traversal.* Then,  $\text{Reset}_r$  is set to false. (Line 47 in Algorithm 1) and, as a consequence, the reset field of the next control token will be false (Line 73 of Algorithm 1). Hence, no resource token will be erased during the next circulation of a control token. If  $PT + \text{SToken}_r < \ell$ , then exactly  $\ell - (PT + \text{SToken}_r)$  are created (see Lines 55 to 58 in Algorithm 1). Hence, the number of resource tokens will be exactly equal to  $\ell$  at the beginning of the next control token circulation. By Lemma 3,  $PT + \text{SToken}_r$  will be equal to  $\ell$  at the end of the next control token circulation, no resource token will be added. Any later circulation of the control token cannot change the number of resource tokens. Hence, the system will contain  $\ell$  resource tokens forever.
- *$PT + \text{SToken}_r > \ell$  at the end of the control token traversal.* Then,  $\text{Reset}_r$  is set to true (Line 47 in Algorithm 1) and, as a consequence, the reset field of the next control token will be true (Line 73 of Algorithm 1). By Lemmas 2 and 3, reducing to the previous case when the circulation of the next control token terminates, and we are done.

□

Following similar reasoning, we can deduce from Lemmas 2, 4, and 5, the following two lemmas:

**Lemma 7** *Starting from any configuration, the system eventually reaches a configuration from which there always exists one priority token.*

**Lemma 8** *Starting from any configuration, the system eventually reaches a configuration after which there always exists one pusher token.*

We now show that once the system contains the correct number of each type of token, the system stabilizes to the  $k$ -out-of- $\ell$  exclusion specification.

**Lemma 9** *Starting from any configuration, every process receives a pusher token infinitely many times.*

**Proof.** By Lemmas 6, 7, and 8, starting from any configuration, the system eventually reaches a configuration  $\gamma$  after which there are always exactly  $\ell$  resource tokens, one priority token, and one pusher token in the network. From  $\gamma$ , the system is then never again reset. So, from  $\gamma$ , the unique pusher token of the system always follows DFS order. Each time a process receives the pusher token, it retransmits it in finite time. Hence, every process receives it infinitely often and the lemma holds. □

**Lemma 10** *Starting from any configuration, every process receives a priority token infinitely many times.*

**Proof.** By Lemmas 6, 7, and 8, starting from any configuration, the system eventually reaches a configuration  $\gamma$  from which there are  $\ell$  resource tokens, one priority token, and one pusher token in the network. From  $\gamma$ , the system is never again reset. So from  $\gamma$ , the unique priority token of the system always follow the DFS order.

By way of contradiction, assume that, from  $\gamma$ , a process eventually stops receiving the priority token. Since the priority token circulates in DFS order and traverses each link in finite time, we can deduce that some other process  $p$  eventually holds it forever. In this case,  $p$  is a requester and its request is never satisfied. Now, by Lemma 9 other process receives the pusher token infinitely often. So, each other process retransmits the resource tokens it holds within finite time, because it eventually either satisfies its request, executes its critical sections, and then releases its tokens, or does not satisfy its request, but receives the pusher, and then releases its resource tokens. Similarly the resource tokens always follows DFS order after  $\gamma$ . Hence,  $p$  receives resource tokens infinitely many times, and never releases the the priority token even if it receives the pusher token. Since  $k \leq \ell$ , the request of  $p$  is eventually satisfied, contradiction.  $\square$

**Lemma 11** *Starting from any configuration, every process receives resource tokens infinitely many times.*

**Proof.** By Lemmas 6, 7, and 8, starting from any configuration, the system eventually reaches a configuration  $\gamma$  from which there are  $\ell$  resource tokens, one priority token, and one pusher token in the network. After  $\gamma$ , the system is then never again reset. Thus, after  $\gamma$ , the resource tokens of the system always follow DFS order.

Assume, by way of contradiction, that some process only receives resource tokens finitely many times. This implies that every resource token is eventually held forever by some process. Consider one process that holds at least one resource token forever. By Lemma 9, that process cannot hold the priority token forever. When it releases the priority token, either its request is satisfied, it executes the critical section within finite time, and then releases its resource tokens, or it is not a requester and thus must release its resource token. Either case is a contradiction, and we are done.  $\square$

**Lemma 12** *Starting from any configuration, the fairness property of the  $k$ -out-of- $\ell$  exclusion specification is eventually satisfied.*

**Proof.** Assume that there a request by some process  $p$  that is never satisfied. By Lemmas 6, 7, and 8, starting from any configuration, the system eventually reaches a configuration  $\gamma$  from which there are always  $\ell$  resource tokens, one priority token, and one pusher token into the network. After  $\gamma$ , the system is then never again reset. Hence, after  $\gamma$ , if  $p$  holds the priority token, it releases it only if its request is satisfied. By Lemma 11,  $p$  eventually receives the priority token. Again by Lemma 11,  $p$  eventually releases the priority token, and so its request must have been satisfied, contradiction.  $\square$

**Lemma 13** *Starting from any configuration, the safety property of the  $k$ -out-of- $\ell$  exclusion specification is eventually satisfied.*

**Proof.** By Lemma 6, there are eventually exactly  $\ell$  resource tokens in the network. Hence, eventually, exactly  $\ell$  resource unit are available in the system.

Finally, any process  $p$  that initially holds some resource tokens eventually releases them because either is is not a requester or it eventually satisfies its request by Lemma 12. Hence, eventually  $p$  sets  $RSet$  to  $\emptyset$  and then  $|RSet| \leq \text{Need}$  forever because each time  $p$  receives a resource token while  $|RSet| \geq \text{Need}$ , it directly retransmits it (see Lines 10 to 19 in Algorithm 1 and Lines 9 to 15 in Algorithm 2). Now,  $\text{Need}$  is always less than or equal to  $k$ . Hence, every process eventually only uses at most  $k$  resource tokens (units) simultaneously.  $\square$

**Lemma 14** *Starting from any configuration, the efficiency property of the  $k$ -out-of- $\ell$  exclusion specification is eventually satisfied.*

**Proof.** We use the definition of efficiency given in [3]. We prove that starting from any configuration,  $(k, \ell)$ -liveness is eventually satisfied.

Consider the configuration  $\gamma$  after which: (1) there are always  $\ell$  resource tokens, one priority token, and one pusher token; and (2) the safety properties of the  $k$ -out-of- $\ell$  exclusion are satisfied (such a configuration exists by Lemmas 6, 7, 8, and 13).

Assume that after  $\gamma$ , the system reaches a configuration  $\gamma'$  after which there is a subset  $I$  of processes such that every process in  $I$  executes its critical section forever (in this case they hold some resource units forever). Let  $\alpha$  be the total number of resource units held forever by the processes in  $I$ .

Assume then that there are some processes not in  $I$  that request some resource units and each of these processes requests at most  $\ell - \alpha$  resource units.

The priority token follows DFS order. Since every process in  $I$  executes the critical section forever, none of these processes keeps the priority token forever (see Lines 92 in Algorithm 1 and 73 in Algorithm 2). Finally, every non-requester directly retransmits the priority token when it receives it (see Line 92 in Algorithm 1 and Line 73 in Algorithm 2). Hence, there is a requesting process  $p$  which is not in  $I$  that eventually receives the priority token. From that point,  $p$  will release it only after its request is satisfied (see Line 92 in Algorithm 1 and Line 73 in Algorithm 2). As a consequence,  $p$  will keep every resource token it receives, even if it receives the pusher token. Checking the proof of Lemma 9, we can see that Lemma 9 still holds even if some processes execute the critical section forever. So, by Lemma 9 every process that is not in  $I \cup \{p\}$  receives the pusher token infinitely often, and so cannot hold resource tokens forever. Finally, every process in  $I$  directly retransmits the resource tokens it receives while it is executing the critical section because they satisfy  $|\text{RSet}| \geq \text{Need}$  by Lemma 13 (see Lines 10 to 19 in Algorithm 1 and Lines 9 to 15 in Algorithm 2). So,  $p$  eventually receives the resource tokens it needs to perform the critical section (remember that  $p$  requests at most  $\ell - \alpha$  resource units) and we are done.  $\square$

From Lemmas 13, 12, and 14, we obtain:

**Theorem 1** *The protocol proposed in Algorithms 1 and 2 is a self-stabilizing  $k$ -out-of- $\ell$  exclusion protocol for tree networks.*

### Waiting Time.

**Theorem 2** *Once the protocol proposed in Algorithms 1 and 2 is stabilized, the waiting time is  $\ell \times (2n - 3)^2$  in the worst case.*

**Proof.** We first show that the waiting time of a requesting process that holds the priority token is  $\ell \times (2n - 3)$  in the worst case. Consider a process  $p$  that requests some resource units and holds the priority token. In the worst case,  $p$  appears only once in the virtual ring defined by the DFS order (if  $p$  is a leaf). Also in the worst case, the  $\ell$  resource tokens may traverse the entire virtual ring before  $p$  receives the tokens it needs. The virtual ring can contain up to  $(2n - 3)$  processes in addition to  $p$ . Any resource token may satisfy one request each time it traverses a process (in the worst case, each process other than  $p$  always requests one token). Hence, the  $\ell$  resource tokens may satisfy up to  $\ell \times (2n - 3)$  requests before  $p$  satisfies its request.

Using similar reasoning, we can see that a requesting process could wait until the priority token traverses the whole virtual ring (up to  $2(n - 2)$  nodes) before it satisfy its request; during that time, up to  $\ell \times (2n - 3)^2$  requests can be satisfied, and we are done.  $\square$

## 5 Conclusion and Perspectives

In this paper, we propose the first (deterministic) self-stabilizing distributed  $k$ -out-of- $\ell$  exclusion protocol for asynchronous oriented tree networks. The proposed protocol uses a realistic model of computation, the message-passing model. The only restriction we make is to assume that the channels initially contain a bounded known number of arbitrary messages. We make this assumption to obtain a solution that uses bounded memory per process (see the results in [7]). However, if we assume unbounded process memory, our solution can be easily adapted to work without assumptions on channels (following the method presented in [9]).

The main interest in dealing with an oriented tree is that solutions on the oriented tree can be directly mapped to solutions for arbitrary rooted networks by composing the protocol with a spanning tree construction (e.g., [1, 4]).

There are several possible extensions of our work. On the theoretical side, one can investigate whether the waiting time of our solution ( $\ell \times (2n - 3)^2$ ) can be improved. Possible extension to networks where processes are subject to other failure patterns, such as process crashes, remains open. On the practical side, our solution is designed in a realistic model and can be extended to arbitrary rooted networks. Hence, implementing our solution in a real network is a future challenge.

## References

- [1] Y Afek and A Bremner. Self-stabilizing unidirectional network algorithms by power supply. *Chicago Journal of Theoretical Computer Science*, 1998:Article 3, 1998.
- [2] A K Datta, R Hadid, and V Villain. A new self-stabilizing k-out-of-l exclusion algorithm on rings. In Shing-Tsaan Huang and Ted Herman, editors, *Self-Stabilizing Systems*, volume 2704 of *Lecture Notes in Computer Science*, pages 113–128. Springer, 2003.
- [3] A K Datta, R Hadid, and V Villain. A self-stabilizing token-based k-out-of-l exclusion algorithm. *Concurrency and Computation: Practice and Experience*, 15(11-12):1069–1091, 2003.
- [4] Sylvie Delaët, Bertrand Ducourthial, and Sébastien Tixeuil. Self-stabilization with r-operators revisited. In Ted Herman and Sébastien Tixeuil, editors, *Self-Stabilizing Systems*, volume 3764 of *Lecture Notes in Computer Science*, pages 68–80. Springer, 2005.
- [5] EW Dijkstra. Self stabilizing systems in spite of distributed control. *Communications of the Association of the Computing Machinery*, 17:643–644, 1974.
- [6] M J Fischer, N A Lynch, J E Burns, and A Borodin. Distributed fifo allocation of identical resources using small shared space. *ACM Trans. Program. Lang. Syst.*, 11(1):90–114, 1989.
- [7] Mohamed G. Gouda and Nicholas J. Multari. Stabilizing communication protocols. *IEEE Trans. Computers*, 40(4):448–458, 1991.
- [8] Rachid Hadid and Vincent Villain. A new efficient tool for the design of self-stabilizing l-exclusion algorithms: The controller. In Ajoy Kumar Datta and Ted Herman, editors, *WSS*, volume 2194 of *Lecture Notes in Computer Science*, pages 136–151. Springer, 2001.
- [9] Shmuel Katz and Kenneth J. Perry. Self-stabilizing extensions for message-passing systems. *Distributed Computing*, 7(1):17–26, 1993.
- [10] Y Manabe, R Baldoni, M Raynal, and S Aoyagi. k-arbiter: A safe and general scheme for h-out of-k mutual exclusion. *Theor. Comput. Sci.*, 193(1-2):97–112, 1998.
- [11] Y Manabe and N Tajima. (k)-arbiter for h-out of-k mutual exclusion problem. In *ICDCS*, pages 216–223, 1999.
- [12] M Raynal. A distributed solution to the k-out of-m resources allocation problem. In F K H A Dehne, F Fiala, and W W Koczkodaj, editors, *ICCI*, volume 497 of *Lecture Notes in Computer Science*, pages 599–609. Springer, 1991.
- [13] Michel Raynal. *Algorithms for Mutual Exclusion*. MIT Press, 1986.
- [14] Michel Raynal. *Algorithmes du parallélisme, le problème de l'exclusion mutuelle*. Dunod informatique, 1990.
- [15] George Varghese. Self-stabilization by counter flushing. *SIAM J. Comput.*, 30(2):486–510, 2000.